

**Manuscript version: Author's Accepted Manuscript**

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

**Persistent WRAP URL:**

<http://wrap.warwick.ac.uk/114413>

**How to cite:**

Please refer to published version for the most recent bibliographic citation information. If a published version is known of, the repository item page linked to above, will contain details on accessing it.

**Copyright and reuse:**

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

**Publisher's statement:**

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: [wrap@warwick.ac.uk](mailto:wrap@warwick.ac.uk).

# Blind Justice: Fairness with Encrypted Sensitive Attributes

Niki Kilbertus<sup>1 2</sup> Adrià Gascón<sup>3 4</sup> Matt Kusner<sup>3 4</sup> Michael Veale<sup>5</sup> Krishna P. Gummadi<sup>6</sup> Adrian Weller<sup>2 3</sup>

## Abstract

Recent work has explored how to train machine learning models which do not discriminate against any subgroup of the population as determined by sensitive attributes such as gender or race. To avoid disparate treatment, sensitive attributes should not be considered. On the other hand, in order to avoid disparate impact, sensitive attributes must be examined—e.g., in order to learn a fair model, or to check if a given model is fair. We introduce methods from secure multi-party computation which allow us to avoid both. By encrypting sensitive attributes, we show how an outcome-based fair model may be learned, checked, or have its outputs verified and held to account, *without users revealing their sensitive attributes*.

## 1. Introduction

Concerns are rising that machine learning systems which make or support important decisions affecting individuals—such as car insurance pricing, résumé filtering or recidivism prediction—might illegally or unfairly discriminate against certain subgroups of the population (Schreurs et al., 2008; Calders & Žliobaitė, 2012; Barocas & Selbst, 2016). The growing field of *fair learning* seeks to formalize relevant requirements, and through altering parts of the algorithmic decision-making pipeline, to detect and mitigate potential discrimination (Friedler et al., 2016).

Most legally-problematic discrimination centers on differences based on *sensitive attributes*, such as gender or race (Barocas & Selbst, 2016). The first type, *disparate treatment* (or *direct discrimination*), occurs if individuals are treated differently according to their sensitive attributes (with all others equal). To avoid disparate treatment, one should not inquire about individuals’ sensitive attributes. While

this has some intuitive appeal and justification (Grgić-Hlača et al., 2018), a significant concern is that sensitive attributes may often be accurately predicted (“reconstructed”) from non-sensitive features (Dwork et al., 2012). This motivates measures to deal with the second type of discrimination.

*Disparate impact* (or *indirect discrimination*) occurs when the *outcomes* of decisions disproportionately benefit or hurt individuals from subgroups with particular sensitive attribute settings without appropriate justification. For example, firms deploying car insurance telematics devices (Handel et al., 2014) build up high dimensional pictures of driving behavior which might easily proxy for sensitive attributes even when they are omitted. Much recent work in fair learning has focused on approaches to avoiding various notions of disparate impact (Feldman et al., 2015; Hardt et al., 2016; Zafar et al., 2017c).

In order to check and enforce such requirements, the modeler must have access to the sensitive attributes for individuals in the training data—however, this may be undesirable for several reasons (Žliobaitė & Custers, 2016). First, individuals are unlikely to want to entrust sensitive attributes to modelers in all application domains. Where applications have clear discriminatory potential, it is understandable that individuals may be wary of providing sensitive attributes to modelers who might exploit them to negative effect, especially with no guarantee that a fair model will indeed be learned and deployed. Even if certain modelers themselves were trusted, the wide provision of sensitive data creates heightened privacy risks in the event of a data breach.

Further, legal barriers may limit collection and processing of sensitive personal data. A timely example is the EU’s General Data Protection Regulation (GDPR), which contains heightened prerequisites for the collection and processing of some sensitive attributes. Unlike other data, modelers cannot justify using sensitive characteristics in fair learning with their “legitimate interests”—and instead will often need explicit, freely given consent (Veale & Edwards, 2018).

One way to address these concerns was recently proposed by Veale & Binns (2017). The idea is to involve a highly trusted third party, and may work well in some cases. However, there are significant potential difficulties: individuals must disclose their sensitive attributes to the third party (even if an individual trusts the party, she may have concerns that

<sup>1</sup>Max Planck Institute for Intelligent Systems <sup>2</sup>University of Cambridge <sup>3</sup>The Alan Turing Institute <sup>4</sup>University of Warwick <sup>5</sup>University College London <sup>6</sup>Max Planck Institute for Software Systems. Correspondence to: Niki Kilbertus <niki.kilbertus@tuebingen.mpg.de>.

the data may somehow be obtained or hacked by others, e.g., [Graham, 2017](#)); and the modeler must disclose their model to the third party, which may be incompatible with their intellectual property or other business concerns.

**Contribution.** We propose an approach to detect and mitigate disparate impact without disclosing readable access to sensitive attributes. This reflects the notion that decisions should be blind to an individual’s status—depicted in courtrooms by a blindfolded Lady Justice holding balanced scales ([Bennett Capers, 2012](#)). We assume the existence of a regulator with fairness aims (such as a data protection authority or anti-discrimination agency). With recent methods from *secure multi-party computation* (MPC), we enable auditable fair learning while ensuring that both individuals’ sensitive attributes and the modeler’s model remain private to all other parties—including the regulator. Desirable fairness and accountability applications we enable include:

1. **Fairness certification.** Given a model and a dataset of individuals, check that the model satisfies a given fairness constraint (we consider several notions from the literature, see Section 2.2); if yes, generate a certificate.
2. **Fair model training.** Given a dataset of individuals, learn a model guaranteed and certified to be fair.
3. **Decision verification.** A malicious modeler might go through fair model training, but then use a different model in practice. To address such accountability concerns ([Kroll et al., 2016](#)), we efficiently provide for an individual to challenge a received outcome, verifying that it matches the outcome from the previously certified model.

We rely on recent theoretical developments in MPC (see Section 3) which we extend to admit linear constraints in order to enforce fairness requirements. These extensions may be of independent interest. We demonstrate the real-world efficacy of our methods, and shall make our code publicly available.

## 2. Fairness and Privacy Requirements

Here we formalize our setup and requirements.

### 2.1. Assumptions and Incentives

We assume three categories of participants: a *modeler*  $M$ , a *regulator*  $REG$ , and *users*  $U_1, \dots, U_n$ . For each user, we consider a vector of sensitive features (or attributes, we use the terms interchangeably)  $\mathbf{z}_i \in \mathcal{Z}$  (e.g., ethnicity or gender) which might be a source of discrimination, and a vector of non-sensitive features  $\mathbf{x}_i \in \mathcal{X}$  (discrete or real). Additionally, each user has a non-sensitive feature  $y_i \in \mathcal{Y}$  which the modeler  $M$  would like to predict—the *label* (e.g., loan default). In line with current work in fair learning, we

assume that all  $\mathbf{z}_i$  and  $y_i$  attributes are binary, though our MPC approach could be extended to multi-label settings. The source of societal concern is that sensitive attributes  $\mathbf{z}_i$  are potentially correlated with  $\mathbf{x}_i$  or  $y_i$ .

Modeler  $M$  wishes to train a model  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ , which accurately maps features  $\mathbf{x}_i$  to labels  $y_i$ , in a supervised fashion. We assume  $M$  needs to keep the model private for intellectual property or other business reasons. The model  $f_\theta$  does not use sensitive information  $\mathbf{z}_i$  as input to prevent disparate treatment (direct discrimination).

For each user  $U_i$ ,  $M$  observes or is provided  $\mathbf{x}_i, y_i$ . The sensitive information in  $\mathbf{z}_i$  is required to ensure  $f_\theta$  meets a given disparate impact fairness condition  $\mathbb{F}$  (see Section 2.2). While each user  $U_i$  wants  $f_\theta$  to meet  $\mathbb{F}$ , they also wish to keep  $\mathbf{z}_i$  private from all other parties. The regulator  $REG$  aims to ensure that  $M$  deploys only models that meet fairness condition  $\mathbb{F}$ . It has no incentive to collude with  $M$  (if collusion were a concern, more sophisticated cryptographic protocols would be required). Further, the modeler  $M$  might be legally obliged to demonstrate to the regulator  $REG$  that their model meets fairness condition  $\mathbb{F}$  before it can be publicly deployed. As part of this,  $REG$  also has a positive duty to enable the training of fair models.

In Section 2.3, we define and address three fundamental problems in our setup: certification, training, and verification. For each problem, we present its functional goal and its privacy requirements. We refer to  $\mathbf{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and  $\mathbf{Z} = \{\mathbf{z}_i\}_{i=1}^n$  as the non-sensitive and sensitive data, respectively. In Section 2.2, we first provide necessary background on various notions of fairness that have been explored in the fair learning literature.

### 2.2. Fairness Criteria

In large part, works that formalize fairness in machine learning do so by balancing a certain condition between groups of people with different sensitive attributes,  $\mathbf{z}$  versus  $\mathbf{z}'$ . Several possible conditions have been proposed. Popular choices include (where  $y \in \{0, 1\}$  and  $\hat{y}$  is the prediction of a machine learning model):

$$P(\hat{y} = y \mid \mathbf{z}) = P(\hat{y} = y \mid \mathbf{z}') \quad (\text{acc}) \quad (1)$$

$$P(\hat{y} = y \mid \mathbf{z}, y = 1) = P(\hat{y} = y \mid \mathbf{z}', y = 1) \quad (\text{TPR}) \quad (2)$$

$$P(\hat{y} = y \mid \mathbf{z}, y = 0) = P(\hat{y} = y \mid \mathbf{z}', y = 0) \quad (\text{TNR}) \quad (3)$$

$$P(\hat{y} = y \mid \mathbf{z}, \hat{y} = 1) = P(\hat{y} = y \mid \mathbf{z}', \hat{y} = 1) \quad (\text{PPV}) \quad (4)$$

$$P(\hat{y} = y \mid \mathbf{z}, \hat{y} = 0) = P(\hat{y} = y \mid \mathbf{z}', \hat{y} = 0) \quad (\text{NPV}) \quad (5)$$

$$P(\hat{y} = 1 \mid \mathbf{z}) = P(\hat{y} = 1 \mid \mathbf{z}') \quad (\text{AR}) \quad (6)$$

Respectively, these consider equality of: (1) accuracy, (2) true positive rate, (3) true negative rate, (4) positive predicted value, (5) negative predicted value, or (6) acceptance rate. Works which use these or related notions include ([Hardt et al., 2016](#); [Zafar et al., 2017c;a;b](#)).

In this work we focus on a variant of eq. (6), formulated as a constrained optimization problem by Zafar et al. (2017c) mimicking the  $p\%$ -rule (Biddle, 2006): for any binary protected attribute  $z \in \{0, 1\}$ , it aims to achieve

$$\min \left\{ \frac{P(\hat{y} = 1 | z = 1)}{P(\hat{y} = 1 | z = 0)}, \frac{P(\hat{y} = 1 | z = 0)}{P(\hat{y} = 1 | z = 1)} \right\} \geq \frac{p}{100}. \quad (7)$$

We believe that in future work, a similar MPC approach could also be used for conditions (1), (2) or (3)—i.e., all the other measures which, to our knowledge, have been addressed with efficient standard (non-private) methods.

### 2.3. Certification, Training, and Verification

**Fairness certification.** Given a notion of fairness  $\mathbb{F}$ , the modeler  $M$  would like to work with the regulator  $REG$  to obtain a certificate that model  $f_\theta$  is fair. To do so, we propose that users send their non-sensitive data  $\mathbf{D}$  to  $REG$ ; and send *encrypted* versions of their sensitive data  $\mathbf{Z}$  to both  $M$  and  $REG$ . Neither  $M$  nor  $REG$  can read the sensitive data. However, we can design a secure protocol between  $M$  and  $REG$  (described in Section 3) to certify if the model is fair. This setup is shown in Figure 1 (Left).

While both  $REG$  and  $M$  learn the outcome of the certification, we require the following *privacy constraints*: (C1) *privacy of sensitive user data*: no one other than  $U_i$  ever learns  $\mathbf{z}_i$  in the clear, (C2) *model secrecy*: only  $M$  learns  $f_\theta$  in the clear, and (C3) *minimal disclosure of  $\mathbf{D}$  to  $REG$* : only  $REG$  learns  $\mathbf{D}$  in the clear.

**Fair model training.** How can a modeler  $M$  learn a fair model without access to users' sensitive data  $\mathbf{Z}$ ? We propose to solve this by having users send their non-sensitive data  $\mathbf{D}$  to  $M$  and to distribute encryptions of their sensitive data to  $M$  and  $REG$  as in certification. We shall describe a secure MPC protocol between  $M$  and  $REG$  to train a fair model  $f_\theta$  privately. This setup is shown in Figure 1 (Center).

*Privacy constraints*: (C1) privacy of sensitive user data, (C2) model secrecy, and (C3) minimal disclosure of  $\mathbf{D}$  to  $M$ .

**Decision verification.** Assume that a malicious  $M$  has had model  $f_\theta$  successfully certified by  $REG$  as above. It then swaps  $f_\theta$  for another potentially unfair model  $f_{\theta'}$  in the real world. When a user receives a decision  $\hat{y}$ , e.g., her mortgage is denied, she can then challenge that decision by asking  $REG$  for a verification  $V$ . The verification involves  $M$  and  $REG$ , and consists of verifying that  $f_{\theta'}(\mathbf{x}) = f_\theta(\mathbf{x})$ , where  $\mathbf{x}$  is the user's non-sensitive data. This ensures that the user would have been subject to the same result with the certified model  $f_\theta$ , even if  $f_{\theta'} \neq f_\theta$  and  $f_{\theta'}$  is not fair. Hence, while there is no simple technical way to prevent a malicious  $M$  from deploying an unfair model, it will get

caught if a user challenges a decision that would differ under  $f_\theta$ . This setup is shown in Figure 1 (Right).

*Privacy constraint*: While  $REG$  and the user learn the outcome of the verification, we require (C1) privacy of sensitive user data, and (C2) model secrecy.

### 2.4. Design Choices

We use a regulator for several reasons. Given fair learning is of most benefit to vulnerable individuals, we do not wish to deter adoption with high individual burdens. While MPC could be carried out without the involvement of a regulator, using all users as parties, this comes at a significantly greater computational cost. With current methods, taking that approach would be unrealistic given the size of the user-base in many domains of concern, and would furthermore require all users to be online simultaneously. Introducing a regulator removes these barriers and leaves users' computational burden at a minimum level, with envisaged applications practical with only their web browsers.

In cases where users are uncomfortable sharing  $\mathbf{D}$  with either  $REG$  or  $M$ , it is trivial to extend all three tasks such that all of  $\mathbf{x}_i, y_i, \mathbf{z}_i$  remain private throughout, with the computation cost increasing only by a factor of 2. This extension would sometimes be desirable as it restricts the view of  $M$  to the final model, prohibiting inferences about  $\mathbf{Z}$  when  $\mathbf{D}$  is known. However, this setup hinders exploratory data analysis by the modeler which might promote robust model-building, and, in the case of verification, validation by the regulator that user-provided data is correct.

## 3. Our Solution

Our proposed solution to these three problems is to use Multi-Party Computation (MPC). Before we describe how it can be applied to fair learning, we first present the basic principles of MPC, as well as its limitations particularly in the context of machine learning applications.

### 3.1. MPC for Machine Learning

Multi-Party Computation protocols allow two parties  $P_1$  and  $P_2$  holding secret values  $x_1$  and  $x_2$  to evaluate an agreed-upon function  $f$ , via  $y = f(x_1, x_2)$  in a way in which the parties (either both or one of them) learn *only*  $y$ . For example, if  $f(x_1, x_2) = \mathbb{I}(x_1 < x_2)$ , then the parties would learn which of their values is bigger, but nothing else.<sup>1</sup> This corresponds to the well-known *Yao's millionaires problem*: two millionaires want to conclude who is richer without disclosing their wealth to each other. The problem was introduced by Andrew Yao in 1982, and kicked off the area of multi-party computation in cryptography.

<sup>1</sup>The function  $\mathbb{I}$  is 1 if its argument is true and 0 otherwise.



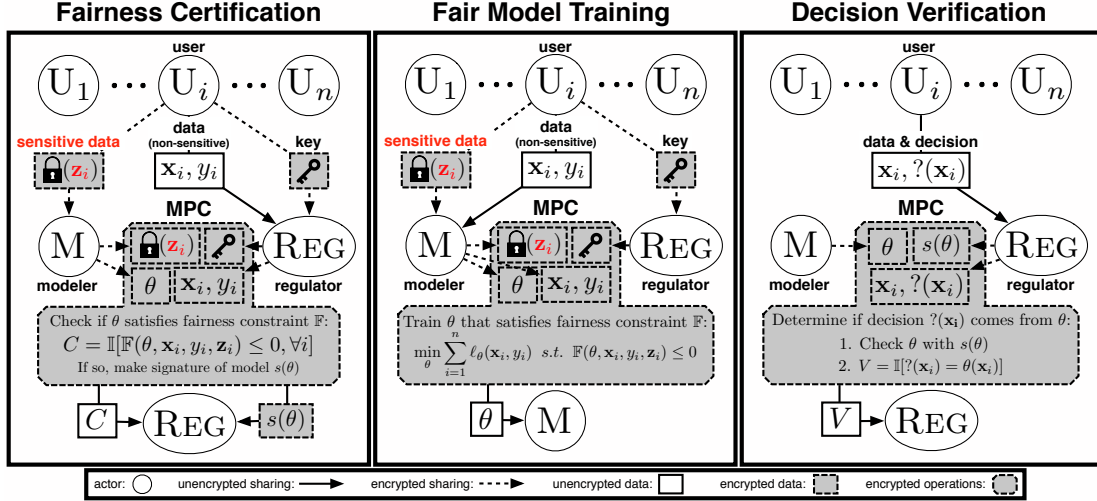


Figure 1. Our setup for Fairness certification (Left), Fair model training (Center), and Decision verification (Right).

In our setting—instead of a simple comparison as in the millionaires problem— $f$  will be either (i) a procedure to check the fairness of a model and certify it, (ii) a machine learning training procedure with fairness constraints, or (iii) a model evaluation to verify a decision. The two parties involved in our computation are the modeler  $M$  and the regulator  $REG$ . The inputs depend on the case (see Figure 1).

As generic solutions do not yet scale to real-world data analysis tasks, one typically has to tailor custom protocols to the desired functionality. This approach has been followed successfully for a variety of machine learning tasks such as logistic and linear regression (Nikolaenko et al., 2013b; Gascón et al., 2017; Mohassel & Zhang, 2017), neural network training (Mohassel & Zhang, 2017) and evaluation (Juvekar et al., 2018; Liu et al., 2017), matrix factorization (Nikolaenko et al., 2013a), and principal component analysis (Al-Rubaie et al., 2017). In the next section we review challenges beyond scalability issues that arise when implementing machine learning algorithms in MPC.

### 3.2. Challenges in Multi-Party Machine Learning

MPC protocols can be classified into two groups depending on whether the target function is represented as either a Boolean or arithmetic circuit. All protocols proceed by having the parties jointly evaluate the circuit, processing it gate by gate while keeping intermediate values hidden from both parties by means of a secret sharing scheme. While representing functions as circuits can be done without losing expressiveness, it means certain operations are impractical. In particular, algorithms that execute different branches depending on the input data will explode in size when implemented as circuits, and in some cases lose their run time guarantees (e.g., consider binary search).

Crucially, this applies to *floating-point arithmetic*. While

this is work in progress, state-of-the-art MPC floating-point arithmetic implementations take more than 15 milliseconds to multiply two 64 bit numbers (Demmler et al., 2015a, Table 4), which is prohibitive for our applications. Hence, machine learning MPC protocols are limited to *fixed-point* arithmetic. Overcoming this limitation is a key challenge for the field. Another necessity for the feasibility of MPC is to approximate non-linear functions such as the sigmoid, ideally by (piecewise) linear functions.

### 3.3. Our MPC Protocols

**Input sharing.** To implement the functionality from Figure 1, we first need a secure procedure for the users to *secret share* a sensitive value, for example her race, with the modeler  $M$  and the regulator  $REG$ . We use *additive secret sharing*. A value  $z$  is represented in a finite domain  $\mathbb{Z}_q$ —we use  $q = 2^{64}$ . To share  $z$ , the user samples a value  $r$  from  $\mathbb{Z}_q$  uniformly at random, and sends  $z - r$  to  $M$  and  $r$  to  $REG$ . While  $z$  can be reconstructed (and subsequently operated on) inside the MPC computation by means of a simple addition, each share on its own does not reveal anything  $z$  (other than that it is in  $\mathbb{Z}_q$ ). One can think of arithmetic sharing as a “distributed one-time pad”.

In Figure 1, we now reinterpret the key held by  $REG$  and the encrypted  $z$  by  $M$ , as their corresponding shares of the sensitive attributes and denote them by  $\langle z \rangle_1$  and  $\langle z \rangle_2$  respectively. The idea of privately outsourcing computation to two non-colluding parties in this way is recurrent in MPC, and often referred to as the two-server model (Mohassel & Zhang, 2017; Gascón et al., 2017; Nikolaenko et al., 2013b; Al-Rubaie et al., 2017).

**Signing and checking a model.** Note that *certification* and *verification* partly correspond to sub-procedures of the *fair training* task: during training we check the fairness

constraint  $\mathbb{F}$ , and repeatedly evaluate partial models on the training dataset (using gradient descent). Hence, *certification* and *verification* do not add technical difficulties over training, which is described in detail in Section 4. However, for verification, we still need to “sign” the model, i.e., REG should obtain a signature  $s(\theta)$  as a result of model certification, see Figure 1 (Left). This signature is used to check in the verification phase, whether a given model  $\theta'$  from M satisfies  $s(\theta') = s(\theta)$  for a certified fair model  $\theta$  (in which case  $\theta = \theta'$  with high probability). Moreover, we need to preserve the secrecy of the model, i.e., REG should not be able to recover  $\theta$  from  $s(\theta)$ . These properties, given that the space of models is large, calls for a cryptographic hash function, such as SHA-256.

Additionally, in our functionality, the hash of  $\theta$  should be computed inside MPC, to hide  $\theta$  from REG. Fortunately, cryptographic hashes such as SHA-256 are a common benchmark functionality in MPC, and their execution is highly optimized. More concretely, the overhead of computing  $s(\theta)$ , which needs to be done both for certification and verification is of the order of fractions of a second (Keller et al., 2013, Figure 14). While cryptographic hash functions have various applications in MPC, we believe the application to machine learning model certification is novel.

Hence, certification is implemented in MPC as a check that  $\theta$  satisfies the criterion  $\mathbb{F}$ , followed by the computation of  $s(\theta)$ . On the other hand, for verification, the MPC protocol first computes the signature of the model provided by M, and then proceeds with a prediction as long as the computed signature matches the one obtained by REG in the verification phase. An alternative solution is possible based on symmetric encryption under a shared key, as highly efficient MPC implementations of block ciphers such as AES are available (Keller et al., 2017).

**Fair training.** To realize the *fair training* functionality from the previous section, we follow closely the techniques recently introduced by Mohassel & Zhang (2017). Specifically, we extend their custom MPC protocol for logistic regression to additionally handle linear constraints. This extension may be of independent interest, and has applications for privacy-preserving machine learning beyond fairness. The concrete technical difficulties in achieving this goal, and how to overcome them, are presented in the next section. The formal privacy guarantees of our fair training protocol are stated in the following proposition.

**Proposition 1.** *For non-colluding M and REG, our protocol implements the fair model training functionality satisfying constraints (C1)-(C3) in Section 2.3 in the presence of a semi-honest adversary.*

The proof holds in the random oracle model, as a standard simulation argument combining several MPC primitives

(Mohassel & Zhang, 2017; Gascón et al., 2017). It leverages security of arithmetic sharing, garbled circuits, and oblivious transfer protocols in the semi-honest model (Goldreich et al., 1987). A general introduction to MPC, as well as a description of the relevant techniques from (Mohassel & Zhang, 2017) used in our protocol, can be found in Section A in the appendix.

## 4. Technical Challenges of Fair Training

We now present our tailored approaches for learning and evaluating fair models with encrypted sensitive attributes. We do this via the following contributions:

- We argue that current optimization techniques for fair learning algorithms are unstable for fixed-point data, which is required by our MPC techniques.
- We describe optimization schemes that are well-suited for learning over fixed-point number representations.
- We combine tricks to approximate non-linear functions with specialized operations to make fixed-point arithmetic feasible and avoid over- and under-flows.

The optimization problem at hand is to learn a classifier  $\theta$  subject to a (often convex) fairness constraint  $\mathbb{F}(\theta)$ :

$$\min_{\theta} \sum_{i=1}^n \ell_{\theta}(\mathbf{x}_i, y_i) \quad \text{subject to } \mathbb{F}(\theta) \leq \mathbf{0}, \quad (8)$$

where  $\ell_{\theta}$  is a loss term (the logistic loss in this work). We collect user data from  $U_1, \dots, U_n$  into matrices  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{Z} \in \{0, 1\}^{n \times p}$  and a label vector  $\mathbf{y} \in \{0, 1\}^n$ .

Zafar et al. (2017c) use a convex approximation of the  $p\%$ -rule, see eq. (7), for linear classifiers to derive the constraint:

$$\mathbb{F}(\theta) = \frac{1}{n} |\hat{\mathbf{Z}}^{\top} \mathbf{X} \theta| - \mathbf{c}, \quad (9)$$

where  $\hat{\mathbf{Z}}$  is the matrix of all  $\hat{\mathbf{z}}_i := \mathbf{z}_i - \bar{\mathbf{z}}$  and  $\mathbf{c} \in \mathbb{R}^d$  is a constant vector corresponding to the tightness of the fairness constraint. Here,  $\bar{\mathbf{z}}$  is the mean of all inputs  $\mathbf{z}_i$ . With  $\mathbf{A} := 1/n \hat{\mathbf{Z}}^{\top} \mathbf{X}$ , the  $p\%$  constraint reads  $\mathbb{F}(\theta) = |\mathbf{A} \theta| - \mathbf{c}$ , where the absolute value is taken element-wise.

### 4.1. Current Techniques

To solve the optimization problem in eq. (8), with the fairness function  $\mathbb{F}$  in eq. (9), Zafar et al. (2017c) use Sequential Least Squares Programming (SLSQP). This technique works by reformulating eq. (8) as a sequence of Quadratic Programs (QPs). After solving each QP, their algorithm uses the Han-Powell method, a quasi-Newton method that iteratively approximates the Hessian  $\mathbf{H}$  of the objective function via the update

$$\mathbf{H}_{t+1} = \mathbf{H}_t + \frac{\mathbf{l}_{\Delta} \mathbf{l}_{\Delta}^{\top}}{\theta_{\Delta}^{\top} \mathbf{l}_{\Delta}} - \frac{\mathbf{H}_t \theta_{\Delta} \theta_{\Delta}^{\top} \mathbf{H}_t}{\theta_{\Delta}^{\top} \mathbf{H}_t \theta_{\Delta}},$$

where  $\mathbf{l}_\Delta = \mathbf{l}(\boldsymbol{\theta}_{t+1}, \boldsymbol{\lambda}_{t+1}) - \mathbf{l}(\boldsymbol{\theta}_t, \boldsymbol{\lambda}_t)$  and  $\mathbf{l}(\boldsymbol{\theta}_t, \boldsymbol{\lambda}_t) = \sum_{i=1}^n \ell_{\boldsymbol{\theta}_t}(\mathbf{x}_i, y_i) + \boldsymbol{\lambda}^\top \mathbb{F}(\boldsymbol{\theta}_t)$  is the Lagrangian of eq. (8). Finally,  $\boldsymbol{\theta}_\Delta = \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t$ .

There are two issues with this approach from an MPC perspective. First, solving a sequence of QPs is prohibitively time-consuming in MPC. Second, while the above Han-Powell update performs well on floating-point data, the two divisions by non-constant, non-integer numbers easily underflow or overflow with fixed-point numbers.

#### 4.2. Fixed-Point-Friendly Optimization Techniques

Instead, to solve the optimization problem in eq. (8), we perform stochastic gradient descent and experiment with the following techniques to incorporate the constraints.

*Lagrangian multipliers.* Here we minimize

$$\mathcal{L} := \frac{1}{n} \sum_{i=1}^n \ell_{\boldsymbol{\theta}}^{\text{BCE}}(\mathbf{x}_i, y_i) + \boldsymbol{\lambda}^\top \max\{\mathbb{F}(\boldsymbol{\theta}), \mathbf{0}\},$$

using stochastic gradient descent, i.e., alternating updates  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} \mathcal{L}$  and  $\boldsymbol{\lambda} \leftarrow \max\{\boldsymbol{\lambda} + \eta_{\boldsymbol{\lambda}} \nabla_{\boldsymbol{\lambda}} \mathcal{L}, \mathbf{0}\}$ , where  $\eta_{\boldsymbol{\theta}}, \eta_{\boldsymbol{\lambda}}$  are the learning rates.

*Projected gradient descent.* For this method, consider specifically the  $p\%$ -rule based notion  $\mathbb{F}(\boldsymbol{\theta}) = |\mathbf{A}\boldsymbol{\theta}| - \mathbf{c}$ . We first define  $\hat{\mathbf{A}}$  as the matrix consisting of the rows of  $\mathbf{A}$  for which  $\mathbb{F}(\boldsymbol{\theta}) > \mathbf{0}$ , i.e., where the constraint is active. In each step, we project the computed gradient of the binary-cross-entropy loss  $\mathcal{L}^{\text{BCE}}$ —either of a single example or averaged over a minibatch—back into the constraint set, i.e.,

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta_{\boldsymbol{\theta}} (\text{Id}_d - \hat{\mathbf{A}}^\top (\hat{\mathbf{A}} \hat{\mathbf{A}}^\top)^{-1} \hat{\mathbf{A}}) \nabla_{\boldsymbol{\theta}} \ell_{\boldsymbol{\theta}}^{\text{BCE}}. \quad (10)$$

*Interior point log barrier (Boyd & Vandenberghe, 2004).* We can approximate eq. (8) for the  $p\%$ -rule constraint  $\mathbb{F}(\boldsymbol{\theta}) = |\mathbf{A}\boldsymbol{\theta}| - \mathbf{c}$  by: minimize  $\sum_{i=1}^n \ell_{\boldsymbol{\theta}}^{\text{BCE}}(\mathbf{x}_i, y_i) - \frac{1}{t} \sum_{j=1}^p (\log(\mathbf{a}_j^\top \boldsymbol{\theta} + c_j) + \log(-\mathbf{a}_j^\top \boldsymbol{\theta} + c_j))$ , where  $\mathbf{a}_j$  is the  $j$ th row of  $\mathbf{A}$ . The parameter  $t$  trades off the approximation of the true objective ( $I_-(u) = 0$  for  $u \leq 0$  and  $I_-(u) = \infty$  for  $u > 0$ ) and the smoothness of the objective function. Throughout training  $t$  is increased, allowing the solution to move closer to the boundary. As the gradient of the objective has a simple closed form representation, we can perform regular (stochastic) gradient descent.

After extensive experiments (see Section 5) we found the Lagrangian multipliers technique to work best, both in yielding high accuracies, reliably staying within the constraints and being robust to hyperparameter changes such as learning rates or the batch size. For a proof of concept, in Section 5 we focus on the  $p\%$ -rule, i.e., eq. (9). Note that the gradients for eq. (2) and eq. (3) take a similarly simple form, i.e., balancing the true positive or true negative rates (corresponding to equal opportunity or equal odds) is simple to implement

for the Lagrangian multiplier technique, but harder for projected gradient descent. However, these fairness notions are more expensive as we have to compute  $\mathbf{Z}^\top \mathbf{X}$  for each update step, instead of pre-computing it once at the beginning of training, see Algorithm 1 in the appendix. We could speed up the computation again by evaluating the constraint only on the current minibatch for each update, in which case we risk violating the fairness constraint.

**MPC-friendliness.** For eq. (9), we can compute the gradient updates in all three methods with elementary linear algebra operations (matrix multiplications) and a single evaluation of the logistic function. While MPC is well suited for linear operations, most nonlinear functions are prohibitively expensive to evaluate in an MPC framework. Hence we tried two piecewise linear approximations for  $\sigma(x)$ . The first was recently suggested for machine learning in an MPC context (Mohassel & Zhang, 2017) and is simply constant 0 and 1 for  $x < -0.5$  and  $x > 0.5$  respectively, and linear in between. The second uses the optimal first order Chebychev polynomial on each interval  $[x, x+1]$  for  $x \in \{-5, -4, \dots, 4\}$ , and is constant 0 or 1 outside of  $[-5, 5]$  (Faiedh et al., 2001). While it is more accurate, we only report results for the simpler first approximation, as it yielded equal or better results in all our experiments.

As the largest number that can be represented in fixed-point format with  $m$  integer and  $m$  fractional bits is roughly  $2^m + 1$ , overflow becomes a common problem. Since we whiten the features  $\mathbf{X}$  column-wise, we need to be careful whenever we add roughly  $2^m$  numbers or more, because we cannot even represent numbers greater than  $2^m$ . In particular, the minibatch size has to be smaller than this limit. For large  $n$ , the multiplication  $\mathbf{Z}^\top \mathbf{X}$  in the fairness function  $\mathbb{F}$  for the  $p\%$ -rule is particularly problematic.

Hence, we split both factors into blocks of size  $b \times b$  with  $b < 2^m$  and normalize the result of each blocked matrix multiplication by  $b$  before adding up the blocks. We then multiply the sum by  $b/n > 2^{-m}$ . As long as  $b, b/n$  (and thus also  $n/b$ ) can be represented with sufficient precision, which is the case in all our experiments, this procedure avoids under- and overflow. Note that we require the sample size  $n$  to be a multiple of  $b$ . In practice, we have to either discard or duplicate part of the data. Since the latter may introduce bias, we recommend subsampling. Once we have (an approximation of)  $\mathbf{A} \in \mathbb{R}^{p \times d}$ , we resort to normal matrix multiplication, as typically  $p, d \lesssim 100$ , see Table 1.

Division is prohibitively expensive in MPC. Hence, we set the minibatch size to a power of two, which allows us to use fast bit shifts for divisions when averaging over minibatches. To exploit the same trick when averaging over/across blocks in the blocked matrix multiplication, we choose  $n$  as the largest possible power of two, see Table 1.

Table 1. Dataset sizes and online timing results of MPC certification and training over 10 epochs with batch size 64.

	Adult	Bank	COMPAS	German	SQF
$n$ training examples	$2^{14}$	$2^{15}$	$2^{12}$	$2^9$	$2^{16}$
$d$ features	51	62	7	24	23
$p$ sensitive attr.	1	1	7	1	1
certification	802 ms	827 ms	288 ms	250 ms	765 ms
training	43 min	51 min	7 min	1 min	111 min

Algorithm 1 in Section B in the appendix describes the computations M and REG have to run for fair model training using the Lagrangian multiplier technique and the  $p\%$ -rule from eq. (9). We implicitly assume all computations are performed jointly on additively shared secrets.

## 5. Experiments

The root cause for most technical difficulties pointed out in the previous section is the necessity to work with fixed-point numbers and the high computational cost of MPC. Hence, major concerns are loss of precision and infeasible running times. In this section, we show how to overcome both doubts and that fair training, certification and verification are feasible for realistic datasets.

### 5.1. Experimental Setup and Datasets

We work with two separate code bases. Our Python code does not implement MPC, to be able to flexibly switch between floating and fixed-point numbers as well as exact non-linear functions and their approximations. We use it mostly for validation and empirical guidance in our design choices. The full MPC protocol is implemented in C++ on top of the Obliv-C garbled circuits framework (Zahur & Evans, 2015a) and the Absentminded Crypto Kit (lib). This is done as described in Section 3 for the Lagrangian multiplier technique (see Section A in the appendix for more details). It accurately mirrors the computations performed by the first implementation on encrypted data.<sup>2</sup> Except for the timing results in Table 1, all comparisons with floating-point numbers or non-linearities were done with the versatile Python implementation. Details about parameters and the algorithm can be found in Section B in the appendix.

We consider 5 real world datasets, namely the adult (*Adult*), German credit (*German*), and bank market (*Bank*) datasets from the UCI machine learning repository (Lichman, 2013), the stop, question and frisk 2012 dataset (*SQF*),<sup>3</sup> and the COMPAS dataset (Angwin et al., 2016) (*COMPAS*). For practical purposes (see Section 4), we subsample  $2^i$  examples from each dataset with the largest possible  $i$ , see Table 1. Moreover, we also run on synthetic data, generated

<sup>2</sup>Code is available at <https://github.com/nikilbertus/blind-justice>

<sup>3</sup><https://perma.cc/6CSM-N7AQ>

as described by Zafar et al. (2017c, Section 4.1), as it allows us to control the correlation between the sensitive attributes and the class labels. It is thus well suited to observe how different optimization techniques handle the fairness-accuracy trade off. For comparison we use the SLSQP approach described in Section 4.1 as a baseline. We run all methods for a range of constraint values in  $[10^{-4}, 10^0]$  and a corresponding range for SLSQP.

In the plots in this section, discontinuations of lines indicate failed experiments. The most common reasons are overflow and underflow for fixed-point numbers, and instability due to exploding gradients. Plots and analyses for the remaining datasets can be found in Section C in the appendix.

### 5.2. Comparing Optimization Techniques

First we evaluate which of the three optimization techniques works best in practice. Figure 2 shows the test set accuracy over the constraint value. By design, the synthetic dataset exhibits a clear trade-off between accuracy and fairness. The *Lagrange* technique closely follows the (dotted) baseline from (Zafar et al., 2017c), whereas *iplb* performs slightly worse (and fails for small  $c$ ). Even though the *projected* gradient method formally satisfies the proxy constraint for the  $p\%$  rule, it does so by merely shrinking the parameter vector  $\theta$ , which is why it also fails for small  $c$ . We analyze this behavior in more detail in Section C in the appendix.

The COMPAS dataset is the most challenging as it contains 7 sensitive attributes, one of which has only 10 positive instances in the training set. Since we enforce the fairness constraint individually for each sensitive attribute (we randomly picked one for visualization), the classifier tends to collapse to negative predictions. All three methods maintain close to optimal accuracy in the unconstrained region, but collapse more quickly than SLSQP. This example shows that the  $p\%$ -rule proxy itself needs careful interpretation when applied to multiple sensitive attributes simultaneously and that our SGD based approach seems particularly prone to collapse in such a scenario. On the Bank dataset accuracy increases for *iplb* and *Lagrange* when the constraint becomes active as  $c$  decreases until they match the baseline. Determining the cause of this—perhaps unintuitive—behavior requires further investigation. We currently suspect the constraint to act as a regularizer. The *projected* gradient method is unreliable on the Bank dataset.

Empirically, the Lagrangian multiplier technique is most robust with maximal deviations of accuracy from SLSQP of  $< 4\%$  across the 6 datasets and all constraint values. We substantiate this claim in Section C of the appendix. For the rest of this section we only report results for Lagrangian multipliers. Figure 2 also shows that using a piecewise linear approximation as described in Section 4 for the logistic function does not spoil performance.



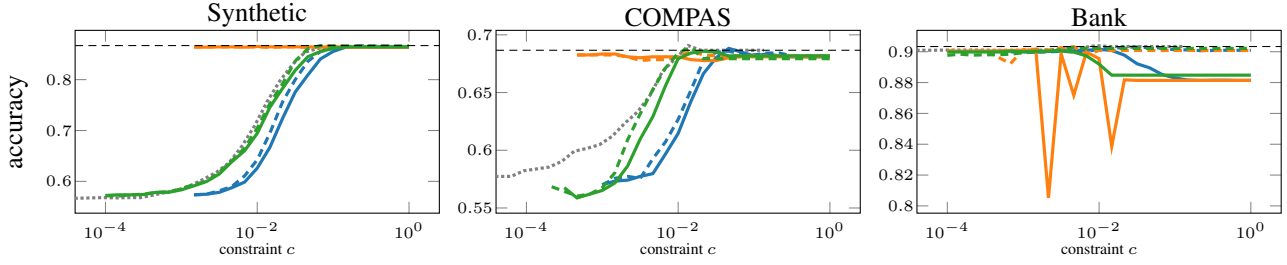


Figure 2. Test set accuracy over the  $p\%$  value for different optimization methods (blue: iplb, orange: projected, green: Lagrange) and either no approximation (*continuous*) or a piecewise linear approximation (*dashed*) of the sigmoid using floating-point numbers. The gray dotted line is the baseline (see Section 4.1) and the black dashed line is unconstrained logistic regression (from scikit-learn).

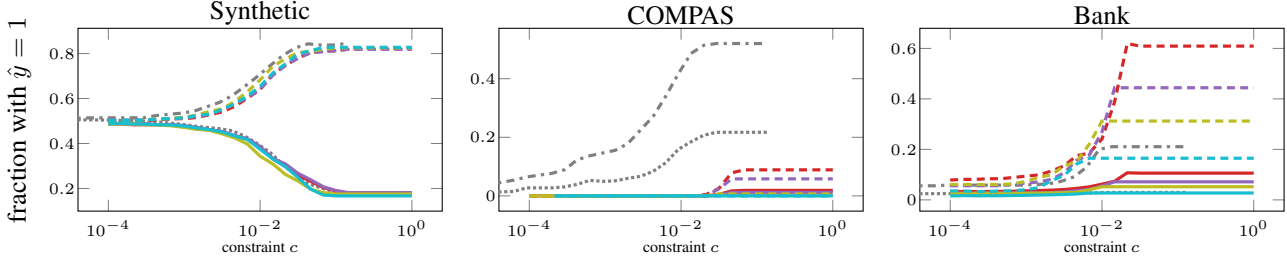


Figure 3. The fraction of people with  $z = 0$  (*continuous/dotted*) and  $z = 1$  (*dashed/dash-dotted*) who get assigned positive outcomes (red: no approx. + float, purple: no approx. + fixed, yellow: pw linear + float, turquoise: pw linear + fixed, gray: baseline).

### 5.3. Fair Training, Certification and Verification

Figure 3 shows how the fractions of users with positive outcomes in the two groups ( $z = 0$  is continuous and  $z = 1$  is dashed) are gradually balanced as we decrease the fairness constraint  $c$ . These plots can be interpreted as the degree to which disparate impact is mitigated as the constraint is tightened. The effect is most pronounced for the synthetic dataset by construction. As discussed above, the collapse for the COMPAS dataset occurs faster than for SLSQP due to the constraints from multiple sensitive attributes. In the Bank dataset, for large  $c$ —before the constraint becomes active—the fractions of positive outcomes for  $z = 1$  differ, which is related to the slightly suboptimal accuracy at large  $c$  that needs further investigation. However, as the constraint becomes active, the fractions are balanced at a similar rate as the baseline. Overall, our Lagrangian multiplier technique with fixed point numbers and piecewise linear approximations of non-linearities robustly manages to satisfy the  $p\%$ -rule proxy at similar rates as the baseline with only minor losses in accuracy on all but the challenging COMPAS dataset.

In Table 1 we show the online running times of 10 training epochs on a laptop computer. While training takes several orders of magnitudes longer than a non-MPC implementation, our approach still remains feasible and realistic. We use the one time offline precomputation of multiplication triples described and timed in Mohassel & Zhang (2017, Table 2). As pointed out in Section 3, certification of a trained model requires checking whether  $\mathbb{F}(\theta) > 0$ . We already perform this check at least once for each gradient

update during training. It only takes a negligible fraction of the computation time, see Table 1. Similarly, the operations required for certification stay well below one second.

**Discussion.** In this section, we have demonstrated the practicability of private and fair model training, certification and verification using MPC as described in Figure 1. Using the methods and tricks introduced in Section 4, we can overcome accuracy as well as over- and underflow concerns due to fixed-point numbers. Offline precomputation combined with a fast C++ implementation yield viable running times for reasonably large datasets on a laptop computer.

## 6. Conclusion

Real world fair learning has suffered from a dilemma: in order to enforce fairness, sensitive attributes must be examined; yet in many situations, users may feel uncomfortable in revealing these attributes, or modelers may be legally restricted in collecting and utilizing them. By introducing recent methods from MPC, and extending them to handle linear constraints as required for various notions of fairness, we have demonstrated that it is practical on real-world datasets to: (i) certify and sign a model as fair; (ii) learn a fair model; and (iii) verify that a fair-certified model has indeed been used; all while maintaining cryptographic privacy of all users’ sensitive attributes. Connecting concerns in privacy, algorithmic fairness and accountability, our proposal empowers regulators to provide better oversight, modelers to develop fair and private models, and users to retain control over data they consider highly sensitive.

## Acknowledgments

The authors would like to thank Chris Russell and Phillipp Schoppmann for useful discussions and help with the implementation, as well as the anonymous reviewers for helpful comments. AG and MK were supported by The Alan Turing Institute under the EPSRC grant EP/N510129/1. MV was supported by EPSRC grant EP/M507970/1. AW acknowledges support from the David MacKay Newton research fellowship at Darwin College, The Alan Turing Institute under EPSRC grant EP/N510129/1 & TU/B/000074, and the Leverhulme Trust via the CFI.

## References

- Absentminded crypto kit. <https://bitbucket.org/jackdoerner/absentminded-crypto-kit>.
- Al-Rubaie, M., Wu, P. Y., Chang, J. M., and Kung, S. Privacy-preserving PCA on horizontally-partitioned data. In *DSC*, pp. 280–287. IEEE, 2017.
- Angwin, J., Larson, J., Mattu, S., and Kirchner, L. Machine bias: There is software used across the country to predict future criminals. and it is biased against blacks. *ProPublica*, May, 23, 2016.
- Barocas, S. and Selbst, A. D. Big data’s disparate impact. *California Law Review*, 104:671–732, 2016.
- Bennett Capers, I. Blind justice. *Yale Journal of Law & Humanities*, 24:179, 2012.
- Biddle, D. *Adverse impact and test validation: A practitioner’s guide to valid and defensible employment testing*. Gower Publishing, Ltd., 2006.
- Boyd, S. and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.
- Calders, T. and Žliobaite, I. Why unbiased computational processes can lead to discriminative decision procedures. In *Discrimination and Privacy in the Information Society*, pp. 43–59. Springer, 2012.
- Damgård, I., Pastro, V., Smart, N. P., and Zakarias, S. Multiparty computation from somewhat homomorphic encryption. In *CRYPTO*, volume 7417 of *Lecture Notes in Computer Science*, pp. 643–662. Springer, 2012.
- Demmler, D., Dessouky, G., Koushanfar, F., Sadeghi, A., Schneider, T., and Zeitouni, S. Automated synthesis of optimized circuits for secure computation. In *ACM Conference on Computer and Communications Security*, pp. 1504–1517. ACM, 2015a.
- Demmler, D., Schneider, T., and Zohner, M. ABY – a framework for efficient mixed-protocol secure two-party computation. In *NDSS*. The Internet Society, 2015b.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference, ITCS ’12*, pp. 214–226. ACM, 2012.
- Faiedh, H., Gafsi, Z., and Besbes, K. Digital hardware implementation of sigmoid function and its derivative for artificial neural networks. *Proceeding of the 13th International Conference on Microelectronics, 2001.*, pp. 189 – 192, 11 2001.
- Feldman, M., Friedler, S., Moeller, J., Scheidegger, C., and Venkatasubramanian, S. Certifying and removing disparate impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 259–268, 2015.
- Fredrikson, M., Jha, S., and Ristenpart, T. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pp. 1322–1333, 2015.
- Friedler, S. A., Scheidegger, C., and Venkatasubramanian, S. On the (im)possibility of fairness. 2016. arXiv:1609.07236v1 [cs.CY].
- Gascón, A., Schoppmann, P., Balle, B., Raykova, M., Doerner, J., Zahur, S., and Evans, D. Privacy-Preserving Distributed Linear Regression on High-Dimensional Data. *Proceedings on Privacy Enhancing Technologies*, 2017 (4):345–364, October 2017.
- Goldreich, O. *The Foundations of Cryptography – Volume 2, Basic Applications*. Cambridge University Press, 2004.
- Goldreich, O., Micali, S., and Wigderson, A. How to play any mental game or A completeness theorem for protocols with honest majority. In *STOC*, pp. 218–229. ACM, 1987.
- Graham, C. NHS cyber attack: Everything you need to know about ‘biggest ransomware’ offensive in history. *Telegraph*, May 20, 2017.
- Grgić-Hlača, N., Zafar, M. B., Gummadi, K. P., and Weller, A. Beyond distributive fairness in algorithmic decision making: Feature selection for procedurally fair learning. In *AAAI*, 2018.
- Handel, P., Skog, I., Wahlstrom, J., Bonawiede, F., Welch, R., Ohlsson, J., and Ohlsson, M. Insurance telematics: Opportunities and challenges with the smartphone solution. *IEEE Intelligent Transportation Systems Magazine*, 6(4):57–70, 2014.
- Hardt, M., Price, E., and Srebro, N. Equality of opportunity in supervised learning. In *NIPS*, 2016.

- Juvekar, C., Vaikuntanathan, V., and Chandrakasan, A. Gazelle: A Low Latency Framework for Secure Neural Network Inference. *IACR Cryptology ePrint Archive*, 2018:73, 2018.
- Keller, M., Scholl, P., and Smart, N. P. An architecture for practical actively secure MPC with dishonest majority. In *ACM Conference on Computer and Communications Security*, pp. 549–560. ACM, 2013.
- Keller, M., Orsini, E., Rotaru, D., Scholl, P., Soria-Vazquez, E., and Vivek, S. Faster secure multi-party computation of AES and DES using lookup tables. In *ACNS*, volume 10355 of *Lecture Notes in Computer Science*, pp. 229–249. Springer, 2017.
- Keller, M., Pastro, V., and Rotaru, D. Overdrive: Making SPDZ great again. In *EUROCRYPT (3)*, volume 10822 of *Lecture Notes in Computer Science*, pp. 158–189. Springer, 2018.
- Kroll, J. A., Huey, J., Barocas, S., Felten, E. W., Reidenberg, J. R., Robinson, D. G., and Yu, H. Accountable algorithms. *University of Pennsylvania Law Review*, 165, 2016.
- Lichman, M. UCI machine learning repository, 2013. URL <http://archive.ics.uci.edu/ml>.
- Lindell, Y. How To Simulate It – A Tutorial on the Simulation Proof Technique. *IACR Cryptology ePrint Archive*, 2016:46, 2016.
- Liu, J., Juuti, M., Lu, Y., and Asokan, N. Oblivious neural network predictions via minion transformations. In *CCS*, pp. 619–631. ACM, 2017.
- Mohassel, P. and Zhang, Y. SecureML: A system for scalable privacy-preserving machine learning. In *IEEE Symposium on Security and Privacy (SP)*, pp. 19–38, 2017.
- Nikolaenko, V., Ioannidis, S., Weinsberg, U., Joye, M., Taft, N., and Boneh, D. Privacy-preserving matrix factorization. In *ACM Conference on Computer and Communications Security*, pp. 801–812. ACM, 2013a.
- Nikolaenko, V., Weinsberg, U., Ioannidis, S., Joye, M., Boneh, D., and Taft, N. Privacy-preserving ridge regression on hundreds of millions of records. In *IEEE Symposium on Security and Privacy*, pp. 334–348. IEEE Computer Society, 2013b.
- Schreurs, W., Hildebrandt, M., Kindt, E., and Vanfleteren, M. *Cogitas, Ergo Sum*. The Role of Data Protection Law and Non-discrimination Law in Group Profiling in the Private Sector. In *Profiling the European Citizen*, pp. 241–270. Springer, 2008.
- Tramèr, F., Zhang, F., Juels, A., Reiter, M. K., and Ristenpart, T. Stealing machine learning models via prediction apis. In *USENIX Security Symposium*, pp. 601–618, 2016.
- Veale, M. and Binns, R. Fairer machine learning in the real world: Mitigating discrimination without collecting sensitive data. *Big Data & Society*, 4(2), 2017.
- Veale, M. and Edwards, L. Clarity, Surprises, and Further Questions in the Article 29 Working Party Draft Guidance on Automated Decision-Making and Profiling. *Computer Law & Security Review*, 2018. doi: 10.1016/j.clsr.2017.12.002.
- Yao, A. C.-C. How to Generate and Exchange Secrets (Extended Abstract). In *FOCS*, pp. 162–167. IEEE Computer Society, 1986.
- Zafar, M. B., Valera, I., Gomez-Rodriguez, M., and Gummadi, K. P. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *WWW*, 2017a.
- Zafar, M. B., Valera, I., Rodriguez, M., Gummadi, K., and Weller, A. From parity to preference-based notions of fairness in classification. In *Advances in Neural Information Processing Systems*, pp. 228–238, 2017b.
- Zafar, M. B., Valera, I., Rodriguez, M. G., and Gummadi, K. P. Fairness Constraints: Mechanisms for Fair Classification. In *AISTATS*, 2017c.
- Zahur, S. and Evans, D. Obliv-c: A language for extensible data-oblivious computation. *IACR Cryptology ePrint Archive*, 2015:1153, 2015a.
- Zahur, S. and Evans, D. Obliv-C: A Language for Extensible Data-Oblivious Computation. *IACR Cryptology ePrint Archive*, 2015:1153, 2015b.
- Žliobaitė, I. and Custers, B. Using sensitive personal data may be necessary for avoiding discrimination in data-driven decision models. *Artificial Intelligence and Law*, 24(2):183–201, 2016.

## A. Details of the MPC Protocols

**Secret sharing.** A secret sharing scheme allows one to split a value  $x$  (the secret) among two parties, so that no party has unilateral access to  $x$ . In our setting, a user Alice will secret share a sensitive value, for example her race, among a modeler  $M$  and a regulator  $REG$ . Several secret sharing schemes exist, including *Shamir secret sharing*, *xor sharing*, *Yao sharing*, or *arithmetic multiplicative/additive sharing*. In this work we alternate between Yao sharing and additive sharing for efficiency. In the latter, the value  $x$  is represented in a finite domain  $\mathbb{Z}_q$  with, for example  $q = 2^{32}$ . To share her race, Alice samples a value  $r$  from  $\mathbb{Z}_q$  uniformly at random, and sends  $x - r$  to  $M$  and  $r$  to  $REG$ . We call each of  $x - r$  and  $r$  a *share*, and denote them as  $\langle x \rangle_1$  and  $\langle x \rangle_2$ . Now  $M$  and  $REG$  can recover  $x$  by adding their shares, but each share on its own does not reveal anything about the value of  $x$  (other than that it is smaller than  $q$ ). Note that the case where  $q = 2$  corresponds to xor sharing.

**Function evaluation.** MPC can be classified in two groups depending on how  $f$  is represented: either as a Boolean or arithmetic circuit. All protocols proceed by having the parties jointly evaluate the circuit, processing it gate by gate. For each gate  $g$  for which the value for the input wires  $x, y$  is shared among the parties, the parties run a subprotocol to produce the value  $z = g(x, y)$  of the output wire, again shared, without revealing any information in the process. In the setting where we use arithmetic additive sharing, the two parties  $M$  and  $REG$  hold shares,  $\langle x \rangle_1, \langle y \rangle_1$  and  $\langle x \rangle_2, \langle y \rangle_2$ , respectively. In this case,  $f$  is represented as an arithmetic circuit, and hence each gate  $g$  in the circuit is either an addition or a multiplication. Note that if  $g$  is an addition gate, then a sharing of  $z = g(x, y)$  can be obtained by having each party simply compute locally, i.e., without any interaction,  $\langle z \rangle_i = \langle x \rangle_i + \langle y \rangle_i$ , for  $i \in \{1, 2\}$ . If  $g$  is a multiplication, the subprotocol to compute shares of  $z$  is much more costly. Fortunately, it can be divided into an offline and an online phase.

**The preprocessing model in MPC.** In this model, two parties  $P_1, P_2$  engage in an offline phase, which is data independent, and compute (and store) *shared multiplication triples* of the form  $(a, b, c)$ , with  $c = ab$ . Here,  $a, b \in \mathbb{F}_q$  are drawn uniformly at random, and each value  $a, b, c$  is shared among the parties as explained above. In the online phase, a multiplication gate  $z = \text{mul}(x, y)$  on shared values  $x, y$  can be evaluated as follows: (1) each  $P_i$  sets  $\langle e \rangle_i = \langle x \rangle_i - \langle a \rangle_i$  and  $\langle f \rangle_i = \langle y \rangle_i - \langle b \rangle_i$ , (2) the parties exchange their shares of  $e$  and  $f$  and reconstruct these values locally, and (3) each  $P_i$  computes  $\langle z \rangle_i = (i - 1)ef + f\langle a \rangle_i + e\langle b \rangle_i + \langle c \rangle_i$ . The correctness of this protocol can be easily checked. Privacy relies on the uniform randomness of  $a, b$ , and hence  $\langle e \rangle_i$  and  $\langle f \rangle_i$  completely mask the values of  $\langle x \rangle_i$  and  $\langle y \rangle_i$ ,

respectively. For a formal proof see (Demmler et al., 2015b).

Hence, for each multiplication in the function to be evaluated, the parties need to jointly generate a multiplication triple in advance. For computations with many multiplications (like in our case) this can be a costly process. However, this constraint is easy to accommodate in our architecture for private fair model training, as  $M$  and  $REG$  can run the offline phase once “overnight”. Arithmetic multiplication via precomputed triples is a common technique, used in several popular MPC frameworks (Demmler et al., 2015b; Damgård et al., 2012). In this setting, several protocols for triple generation (which we did not describe) are available (Keller et al., 2018), and under continuous improvement. These protocols are often based on either Oblivious Transfer or Homomorphic Encryption.

**The two-server model for multi-party learning.** Due to a sequence of theoretical and engineering breakthroughs, in the last three decades MPC has gone from being a mathematical curiosity to a technology of practical interest with commercial applications. Several generic protocols for MPC exist, such as the ones based on arithmetic sharing (Damgård et al., 2012), garbled circuits (Yao, 1986), or GMW (Goldreich et al., 1987), with several available implementations (Demmler et al., 2015b; Zahur & Evans, 2015b). These protocols have different trade-offs in terms of the number of parties they support, network requirements, and scalability for different kinds of computations. In our work, we focus on the 2-party case, as the MPC computation is done by  $M$  and  $REG$ . The idea of privately outsourcing computation to two non-colluding parties in this way is recurrent in MPC, and often referred to as the two-server model (Mohassel & Zhang, 2017; Gascón et al., 2017; Nikolaenko et al., 2013b; Al-Rubaie et al., 2017).

While generic protocols exist, these do not yet scale to input sizes typically encountered in machine learning applications like ours. To circumvent this limitation, techniques tailored to specific applications have been proposed. Our protocols fall in this category, extending the SGD protocol from (Mohassel & Zhang, 2017), in which the following useful accelerating techniques are presented.

- **Efficient rescaling:** As our arithmetic shares represent fixed-point numbers, we need to rescale by the precision  $p$  after every multiplication. This involves dividing by  $2^p$ , an expensive operation to do in MPC, and in particular in arithmetic sharing. Mohassel et al. show an elegant solution to this problem: the parties can rescale locally by dropping  $p$  bits of their shares. It is not hard to see that this might produce the wrong result. However, the parameters of the arithmetic secret sharing scheme can be set such that with a tunable arbitrarily large probability the error is at most  $\pm 1$ . This



trick can be used for any division by a power of two.

- **Alternating sharing types:** As already pointed out in previous work (Demmler et al., 2015b), alternating between secret sharing schemes can provide significant acceleration for some applications. Intuitively, arithmetic operations are fast in arithmetic shares, while comparisons are fast in schemes that represent functions as Boolean circuits. Examples of the latter are the GMW protocol and Yao’s garbled circuits. In our implementation, we follow this recipe and implement matrix-vector multiplication using arithmetic sharing, while for evaluating our variant of sigmoid, we rely on the protocol from (Mohassel & Zhang, 2017) implemented with garbled circuits using the Obliv-C framework (Zahur & Evans, 2015b).
- **Matrix multiplication triples:** Another observation made by Mohassel et al. is that the idea described above for preprocessing multiplications over arithmetic shares can be reinterpreted at the level of matrices. This results in a faster online and offline phase (see (Mohassel & Zhang, 2017) for details).

**How to prove that a protocol is secure.** We did not provide a formal definition of security in this paper, and instead referred the reader to (Mohassel & Zhang, 2017). In MPC, privacy in the case of semi-honest adversaries is argued in the simulation paradigm (see (Goldreich, 2004) or (Lindell, 2016) for formal definitions and detailed proofs). Intuitively, in this paradigm one proves that every inference that a party—in our case either REG or M—could draw from observing the execution trace of the protocol could also be drawn from the output of the execution and the party’s input. This is done by proving the existence of a *simulator* that can produce an execution trace that is indistinguishable from the actual execution trace of the protocol. A crucial point is that the simulator only has access to the input and output of the party being simulated.

## B. Details of Fair Model Training

### B.1. The Fair Training Algorithm

Algorithm 1 describes the computations M and REG have to perform for fair model training using the Lagrangian multiplier technique and the  $p\%$ -rule from eq. (9). In the next subsection we describe the parameter values. We implicitly assume all computations are performed jointly on additively shared secrets by M and REG as described in Section 3. This means that M and REG each receive a secret share of the protected attributes  $\mathbf{Z}$ . Following the protocols outlined in Section 3, they can then jointly evaluate the steps in Algorithm 1. This allows them to operate on the sensitive values within the MPC computation, while preventing unilateral

access to them by M and REG. The result of these computations is the same as evaluating the algorithm as described with data in the clear.

BLOCKEDMULTSHIFTAVG stands for the blocked matrix multiplication to avoid overflow for fixed-point numbers described towards the end of Section 4. Note that it already contains the division by  $n$ . The averaging within the blocked matrix multiplications as well as over the results thereof are done by fast bit shifts instead of slow MPC division circuits. This is possible, because we chose all parameters such that divisions are always by powers of two.

We found the piecewise linear approximation of the sigmoid function introduced in (Mohassel & Zhang, 2017)

$$\text{SIGMOIDAPPROX}(x) := \begin{cases} 0 & \text{if } x \leq -\frac{1}{2}, \\ x + \frac{1}{2} & \text{if } -\frac{1}{2} < x < \frac{1}{2}, \\ 1 & \text{if } x \geq \frac{1}{2}. \end{cases}$$

to work best, see Figure 4.

**Algorithm 1** Fair model training with private sensitive values using Lagrangian multipliers for  $\mathbb{F}(\theta) = 1/n|\mathbf{Z}^\top \mathbf{X}| - \mathbf{c}$ .

**Parties:** M, REG.

**Input:** (M)  $\langle \mathbf{Z} \rangle_1 \in \mathbb{Z}_q^{n \times p}$

**Input:** (REG)  $\mathbf{X} \in \mathbb{Z}_q^{n \times d}$ ,  $\mathbf{y} \in \mathbb{Z}_q^n$ ,  $\langle \mathbf{Z} \rangle_2 \in \mathbb{Z}_q^{n \times p}$

**Input:** (Public) Learning rates  $\eta_\theta, \eta_\lambda$ , number of training examples  $n$ , minibatch size  $2^s$ , constraints  $\mathbf{c} \in \mathbb{Z}_q^p$ , number of epochs  $N_e$ .

```

1:  $\theta \leftarrow \mathbf{0}, \lambda \leftarrow \mathbf{0}$ 
2:  $\mathbf{A} \leftarrow \text{BLOCKEDMULTSHIFTAVG}(\mathbf{Z}^\top, \mathbf{X})$ 
3: for all  $j$  from 1 to  $N_e$  do
4:   for all  $i$  from 1 to  $n/2^s$  do
5:      $(\mathbf{X}_i, \mathbf{y}_i) \leftarrow \text{SAMPLEMINIBATCH}(\mathbf{X}, \mathbf{y})$ 
6:      $\mathbb{F} \leftarrow |\mathbf{A}\theta| - \mathbf{c}$ 
7:      $\nabla_\lambda \leftarrow \max\{\mathbb{F}, \mathbf{0}\}$ 
8:      $\sigma \leftarrow \text{SIGMOIDAPPROX}(\mathbf{X}_i\theta)$ 
9:      $\nabla_\theta^{\text{BCE}} \leftarrow \text{SHIFTDIVIDE}(\mathbf{X}_i^\top(\sigma - \mathbf{y}_i), 2^s)$ 
10:     $\nabla_\theta^{\text{CON}} \leftarrow \begin{cases} \mathbf{A}^\top \lambda, & \text{if } \mathbf{A} > \mathbf{0} \wedge \mathbb{F} > 0 \\ -\mathbf{A}^\top \lambda, & \text{if } \mathbf{A} < \mathbf{0} \wedge \mathbb{F} > 0 \\ \mathbf{0}, & \text{if } \mathbb{F} \leq 0 \end{cases}$ 
11:     $\theta \leftarrow \theta - \eta_\theta(\xi_j^{\text{BCE}} \nabla_\theta^{\text{BCE}} + \xi_j^{\text{CON}} \nabla_\theta^{\text{CON}})$ 
12:     $\lambda \leftarrow \max\{\lambda + \eta_\lambda \nabla_\lambda, \mathbf{0}\}$ 
13:   end for
14: end for

```

**Output:** Parameters  $\theta$

### B.2. Description of Training Parameters

All our experiments use a batch size of 64, a fixed number of epochs scaling inversely with dataset size  $n$  (such that we always perform roughly 15 000 gradient updates), fixed learning rates of  $\eta_\theta = 10^{-4}$ ,  $\eta_\lambda = 0.05$ , and an annealing

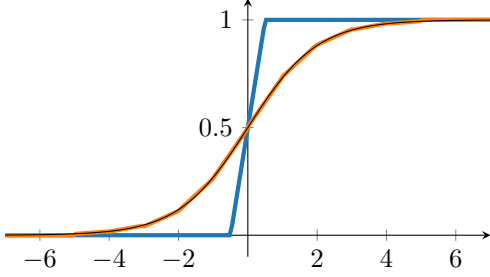


Figure 4. Piecewise linear approximations for the non-linear sigmoid function (in black) from Mohassel & Zhang (2017) in blue and from Faiedh et al. (2001) in orange.

schedule for  $1/t$  in the interior point logarithmic barrier method as described by Boyd & Vandenberghe (2004). The weights for the gradients of the regular binary cross entropy loss (BCE) and the loss from the constraint terms (CON) follow the schedules

$$\xi_j^{\text{BCE}} = \frac{N_e}{N_e + j}, \quad \xi_j^{\text{CON}} = \frac{N_e + 10j}{N_e}.$$

Weight decay, adaptive learning rate schedules, and momentum neither consistently improved nor impaired training. Therefore, all reported numbers were achieved with vanilla SGD, for fixed learning rates, and without any regularization. After extensive testing on all datasets, we converged to a fixed-point representation with 16 bits for the integer and fractional part respectively. The smaller the number of bits, the faster the MPC implementation and the higher the risk of loss of precision or over- and underflows. We found 16 bits to be the minimally needed precision for all our experiments to work.

## C. Additional Experimental Results

### C.1. Results on Remaining Datasets

Analogously to Figures 2 and Figure 3 we report the results on test accuracy as well as the mitigation of disparate impact for the Lagrangian multiplier method in Figure 5. In the Adult dataset we are able to mitigate disparate impact with slightly worse accuracy as compared to the baseline. Note that the German dataset contains only 512 training and 200 test examples, which explains the discrete jumps in accuracy in minimal steps of  $1/200 = 0.005$ . Hence, even though the Lagrangian multiplier technique here consistently removes disparate impact to a similar extent as the baseline, interpretations of results on such small datasets require great care. For the much larger stop, question and frisk dataset we again observe the curious initial increase in accuracy similar to our observations for the Bank dataset. In this dataset about 93% of all samples have positive labels, which explains the near optimal accuracy when collapsing to always predict 1,

which happens for the baseline as well for our method at a similar rate as  $c$  decreases.

### C.2. Disadvantages of Other Optimization Methods

In Section 5 we suggest the Lagrangian multiplier technique for fair model training using fixed-point numbers. Here we substantiate this suggestion with further empirical evidence. Figure 6 shows analogous results to Figure 3 and the second row of Figure 5. These plots reveal the shortcomings of the interior point logarithmic barrier and the projected gradient methods.

**Interior Point Logarithmic Barrier method.** While the interior point logarithmic barrier method does balance the fractions of people being assigned positive outcomes between the two different demographic groups when the constraint is tightened, it soon breaks down entirely due to overflow and underflow errors. The number of failed runs was substantially higher than for the Lagrangian multiplier technique. As explained in (Boyd & Vandenberghe, 2004), when we increase the parameter  $t$  of the interior point logarithmic barrier method during training, the barrier becomes steeper, approaching the function

$$I_-(x) = \begin{cases} 0 & \text{for } x \leq 0, \\ \infty & \text{for } x > 0. \end{cases}$$

From this it becomes obvious that when facing tight constraints, the gradients might change from almost zero to extremely large values within a single update of the parameters  $\theta$ . Moreover, iplb requires careful tuning and scheduling of  $t$ . Hence, the interior point logarithmic barrier method, while achieving good results over some domains, is not well suited for MPC.

**Projected gradient method.** In Figure 6, we observe that the projected gradient method seems to fail in most cases, since it does not actually balance the fractions of positive outcomes across the sensitive groups. There is a simple explanation why it can satisfy the constraint  $\mathbb{F}(\theta) \leq 0$  for the  $p\%$ -rule even with small  $c$  and still retain near optimal accuracy. Note that the accuracy only depends on the direction of  $\theta$ , i.e., it is invariant to arbitrary rescaling of  $\theta$ . Since the constraint  $\mathbb{F}(\theta) = |\mathbf{A}\theta| - c \leq 0$  is always satisfied for  $\theta = 0$ , dividing any  $\theta$  by a large enough factor will result in a classifier that achieves equal accuracy and satisfies the constraint (by continuity). However, minimizing the loss in the original logistic regression optimization problem (or equivalently maximizing the likelihood), which is not invariant under rescaling of  $\theta$ , counteracts shrinking  $\theta$  as it enforces high confidence of decisions, i.e., large  $\theta$ . The projection method produces high accuracy classifiers with small weights that formally fulfill the fairness constraint, but do not properly mitigate disparate impact as measured

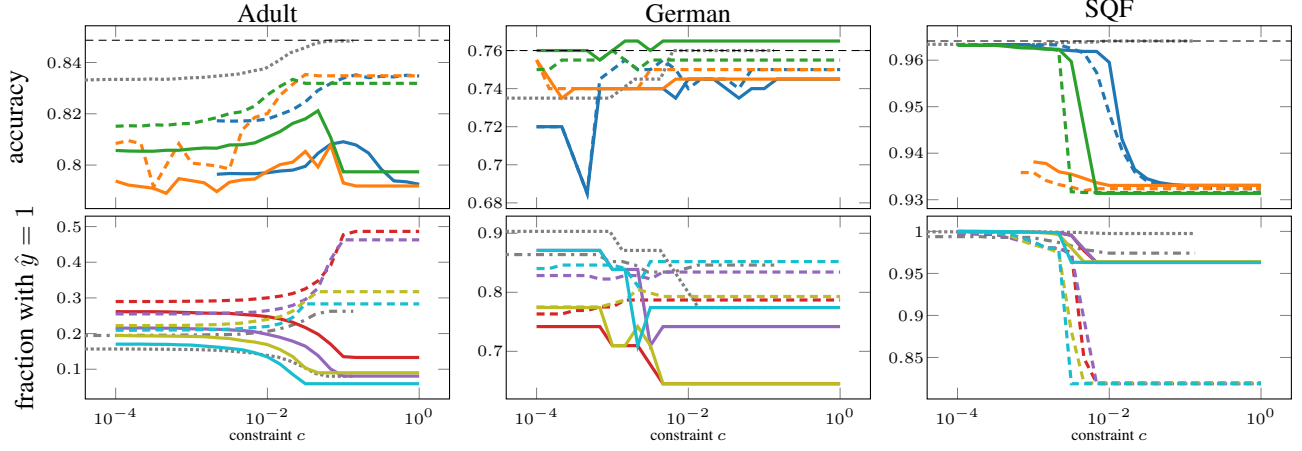


Figure 5. **First row:** The color code is blue: iplb, orange: projected, green: Lagrange with *continuous* lines for no approximation and *dashed* lines for piecewise linear approximation. The gray dotted line is the baseline and the dashed black line marks unconstrained logistic regression. **Second row:** *Continuous/dotted* lines correspond to  $z = 0$  and *dashed/dash-dotted* lines to  $z = 1$ . The color code is (red: no approx. + float, purple: no approx. + fixed, yellow: pw linear + float, turquoise: pw linear + fixed, gray: baseline).

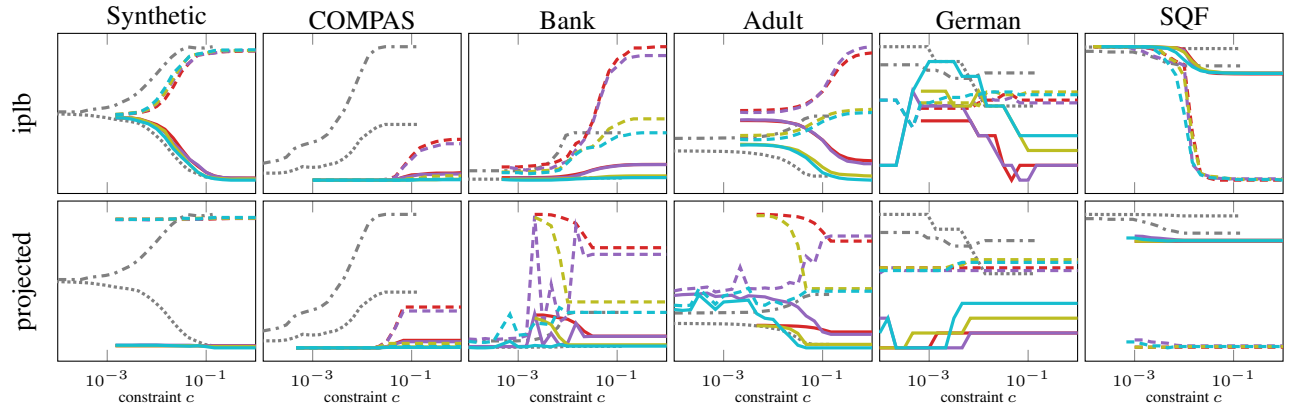


Figure 6. We plot the fraction of people with  $z = 0$  (*continuous/dotted*) and with  $z = 1$  (*dashed/dash-dotted*) who get assigned positive outcomes over the constraint  $c$  for 5 different datasets. The different colors correspond to (red: no approximation + floats, purple: no approximation + fixed-point, yellow: piecewise linear + floats, turquoise: piecewise linear + fixed-point, gray: baseline).

by the true  $p\%$ -rule instead of the computational proxy. It also often fails for small constraint values, as the projection matrix in eq. (10) turns out to become near singular producing over- and underflow errors.

#### **D. Clarification of Privacy or Secrecy**

In this work, privacy or secrecy constraints are separate from other theorized, setup-dependent attacks, e.g., model extraction ([Tramèr et al., 2016](#)) or inversion ([Fredrikson et al., 2015](#)). If relevant, modelers may need to consider these separately.